

# Chain of Second Thoughts: Augmenting Chain-of-Thought Prompting with General Purpose Verifiers

Anonymous ACL submission

## Abstract

Many of the recent capabilities demonstrated by Large Language Models (LLMs) arise primarily from their ability to exploit contextual information. In this paper, we explore ways to improve reasoning capabilities of LLMs through (1) exploration of different chains of thought and (2) validation of the individual steps of the reasoning process. We propose three general principles that a model should adhere to while reasoning: (i) Relevance, (ii) Mathematical Accuracy, and (iii) Logical Consistency. We apply these constraints to the reasoning steps generated by the LLM to improve the accuracy of the final generation. The constraints are applied in the form of verifiers: the model itself is asked to verify if the generated steps satisfy each constraint. To further steer the generations towards high-quality solutions, we use the perplexity of the reasoning steps as an additional verifier. We evaluate our method on 4 distinct types of reasoning tasks, spanning a total of 9 different datasets. Experiments show that our method is always better than vanilla generation, and, in 6 out of the 9 datasets, it is better than best-of N sampling which samples N reasoning chains and picks the lowest perplexity generation.

## 1 Introduction

Large Language Models (LLMs) have demonstrated impressive capabilities of performing a diverse range of tasks by framing them as text generation (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023; OpenAI, 2023; Bubeck et al., 2023, *inter alia*). Chain-of-Thought prompting (Nye et al., 2021; Wei et al., 2022; Chowdhery et al., 2022) further improved their performance on challenging reasoning tasks using a simple trick of generating intermediate steps before giving the final answer allowing the LLM to spread computation over more tokens (Goyal et al., 2023). However, this approach lacks a mechanism to rectify

errors in reasoning. While LLMs may eventually reach the correct answer, they might do so via incorrect intermediate reasoning steps, or worse, never reach the correct answer due to earlier mistakes (Turpin et al., 2023). To illustrate this, we provide a concrete example in Figure 1, where the final answer is correct, but the intermediate steps are (i) irrelevant (Shi et al., 2023), (ii) contradicting previous steps (Mündler et al., 2023), and (iii) with mathematical errors (Patel et al., 2021). Recent work (Yao et al., 2023; Xie et al., 2023; Pan et al., 2023) has attempted to alleviate these problems by employing a search mechanism or a self-correction mechanism in the spirit of "System 2" thinking. Other directions include training a dataset-specific verifier to improve the performance when aggregating multiple reasoning chains (Li et al., 2023). However, all these approaches have dataset-specific adaptations and don't generalize out-of-the-box.

In this work, we explore if catching early mistakes in reasoning chains through problem-agnostic verification can improve reasoning in LLMs. We propose three general principles that a model should adhere to while reasoning: (i) Relevance, (ii) Mathematical Accuracy, and (iii) Logical Consistency and use models, called verifiers, to test for each principle. Each verifier operates on a step (Uesato et al., 2022) generated from the step-by-step manner of Chain-of-Thought prompting and assigns a score to that step. We design the verifiers to operate at this granularity so they can detect intermediate mistakes and discourage the LLM from committing to an erroneous reasoning chain. To further steer the generation towards better steps, we use the perplexity of the reasoning step as an additional verifier. We then explore various ways, including Self-Consistency (Wang et al., 2022), to aggregate verifier scores and report their downstream task performance.

We make the following contributions: (i) we propose a general framework for guiding reasoning in

**Question:** A concert ticket costs \$40. Mr. Doe bought 10. Mr. Benson bought 12 tickets and received a 5% discount for every ticket bought that exceeds 10. How much did Mr. Benson pay in all?

**Step by Step Solution:**

1. Mr. Doe paid in total  $10 * \$40 = \$400$ . (Irrelevant)
2. Without the discount, Mr. Benson paid in total  $12 * 40 = \$480$ .
3. Mr. Benson received a discount of  $2 * (\$40 - \$38) = \$4$  (\$2 per ticket)
4. So Mr. Benson received a discount of \$0.02 per ticket. (Inconsistent)
5. With discount, Mr. Benson paid  $\$480 - 2 * \$2 = \$478$ . (Mathematically inaccurate)
6. Therefore, Mr. Benson paid \$476.

Figure 1: An illustrative example highlighting how the final answer can ultimately be correct (i.e.  $12 * \$40 - (12 - 10) * \$40 * 0.05 = \$476$ ), but it is reached through steps that are (i) irrelevant (Step 1), (ii) contradicting previous steps (Step 4 contradicts Step 3), or (iii) with mathematical errors (Step 5).

LLMs using verifiers which offers the flexibility to use a problem-agnostic implementation across any reasoning task but also offers the adaptability to use task- and dataset-specific implementations, and (ii) we show how using our proposed verifiers can improve reasoning outcomes in LLMs and can also improve existing ensembling techniques like Self-Consistency. Importantly, our work is not intended to be an exploration on the best way to use a computational budget to achieve a desired performance, but an exploration of whether the LLM are capable (even if inefficiently) of detecting their own mistakes together with a simple recovering mechanism.

## 2 Related Work

We focus here only on LLM-based approaches, and divide previous related work according to (i) the generalizability of the prompts used, and (ii) how the final answer is generated.

**Types of Prompts** In prior work, the prompts used can be categorized based on their level of generality. Some approaches utilize a singular prompt, applying it uniformly across a wide spectrum of datasets and tasks. Wei et al. (2022) proposed chain-of-thought prompting with in-context examples. Kojima et al. (2022) then explored zero-shot prompts capable of exhibiting similar behaviors. Other recent works explore using LLMs to self-evaluate (Yin et al., 2023) and potentially improve upon their generation with the resulting feedback (Saunders et al., 2022; Chen et al., 2023; Pan et al., 2023; Shinn et al., 2023). Bai et al. (2022) use

an LLM with in-context examples to detect and edit the responses of a chat model that are harmful or toxic. Madaan et al. (2023) proposes a framework to iteratively self-improve the generations of a LLM. Yao et al. (2023) tightly integrates an LLM with custom dataset-specific prompts to act as a guiding mechanism in the underlying search space. Hao et al. (2023) expands on this by using a Monte Carlo tree search strategy. Other recent work questioned the extent to which using an LLM to evaluate and improve its own generations is viable (Huang et al., 2023), a conclusion which we observed as well in our preliminary work and sidestepped by re-sampling instead of asking the LLM to refine.

Importantly, previous work explored self-evaluation through the lens of task-specific evaluation and prompts, a direction that inherently constrains the broader utility of Large Language Models (LLMs) as general-purpose reasoners. On the other hand, our approach follows a distinct trajectory: we deliberately eschew the use of prompts tailored to individual datasets or tasks.

**How the Final Answer is Generated** A second dimension is that of how the model arrives at the final solution, where we distinguish between methods that take a linear approach (Wei et al., 2022; Kojima et al., 2022; Goyal et al., 2023) from the methods that do not (Yao et al., 2022; Wang et al., 2023; Long, 2023). By linear approaches, we refer to those methods where the final answer is generated token-by-token in one go. On the other hand, non-linear approaches typically include a search mechanism (Xie et al., 2023; Yao et al., 2023; Besta et al., 2023) or a self-reflection process (Madaan et al., 2023; Pan et al., 2023).

For example, recent work explored tightly integrating the LLM to act as a guiding mechanism in the underlying search space (Xie et al., 2023; Yao et al., 2023). This involves one LLM generating candidate steps while another LLM assigns single float value as a value score. This value score is derived from an LLM with a dataset-specific prompt and in-context examples.

Within this dimension, our approach aligns with the non-linear paradigm. We leverage verifiers to evaluate each step in the solution-generation process, with the overarching aim of guiding the generation towards solutions that receive high scores, as determined by the verifiers. Differently from previous work on self-evaluation, we explore a setting of self-evaluation that is problem-agnostic.

### 3 Proposed Method

We propose a novel approach that seamlessly integrates with any given Large Language Model’s (LLM) solution generation process, at both step generation and step evaluation (Wei et al., 2022; Wang et al., 2022; Kojima et al., 2022). Our approach consists of two components: a solution generator  $G$  and a set of verifiers  $\mathcal{V}$ , where each verifier specializes in a particular qualitative aspect of reasoning. The solution generator  $G$  is responsible for generating candidate steps, and each verifier  $v \in \mathcal{V}$  is responsible for checking whether the candidate step is in compliance with the specific reasoning property. We explore properties that are generally applicable to a wide range of reasoning tasks. We provide an illustrative example in Figure 2.

This section is further organized as follows. We describe the notations and abstractions used in Section 3.1, the solution generation component in Section 3.2, the proposed verifiers in Section 3.3, the procedure to obtain fine-grained scores for reasoning chains in Section 3.4 and how we use verifiers in aggregate in Section 3.5.

#### 3.1 Notation

In the following, we define the notation we adopt throughout the paper.

A token as  $t \in \text{Vocab}$  where  $\text{Vocab}$  represents the set of possible tokens defined by a given vocabulary. We use  $[t_1, \dots, t_n] \in T$  to represent a given text, with  $T$  denoting the set of all potential texts of varying lengths. Under this notation, we represent a problem as  $q \in T$  and a reasoning step as  $r \in T$ , both presented in free-text form. A solution generator  $G$  in the form of a function that takes text as input and returns text as output:  $G : T \rightarrow T$ . We interpret the output text as a sequence of reasoning steps  $R^q = [r_1^q, \dots, r_n^q]$ . For simplicity, we define a reasoning step as the sequence of tokens until a new line, similar to Uesato et al. (2022). A verifier  $v \in \mathcal{V}$  implemented a function that takes text as input and returns an indicator:  $v : T \rightarrow \{0, 1\}$ . The returned value represents whether the reasoning step satisfies the verifier constraint (1) or not (0).

We will next describe the Solution Generator and Step Verification components.

#### 3.2 Solution Generation

The solution generator (typically an LLM) operates over a prompt  $q$  and generates a sequence of tokens as output:  $G : T \rightarrow T$ . For our purpose,

we concentrate on the correctness at the level of a reasoning step, instead of individual tokens. A reasoning step, as defined in this work, represents the sequence of tokens up to the occurrence of a new line (Uesato et al., 2022). The next reasoning step can then be generated by conditioning the generator  $G$  on both the question  $q$  and on the sequence of previously generated reasoning steps  $R_{1:i}$ , which is initially empty. This conditioning can be expressed as follows:  $P(r_{i+1}^q | q, R_{1:i}^q)$ .

For the solution generation process, we adopt the zero-shot prompt used in Kojima et al. (2022) which simply appends "Let’s think step by step" to the problem question to elicit a chain-of-thought like behavior in the model’s response without any annotated exemplars (Wei et al., 2022). Nevertheless, our proposed method is agnostic to the specific implementation of the solution generator.

#### 3.3 Step Verification

Our research aims to investigate whether specifying a subset of conditions that an ideal reasoning chain should satisfy and then employing these conditions to score the corresponding reasoning chain can result in improved performance on downstream reasoning tasks.

To this end, we explore three general and necessary (but not sufficient) conditions that a given step should satisfy in order for the resulting solution to be sound from a reasoning perspective: (i) Relevance, (ii) Mathematical Accuracy (if applicable), and (iii) Logical Coherence. In this work, we use for our verifiers a set of LLMs provided with a detailed instruction of the task and constraints.<sup>1</sup> We then map the output of the LLM to  $\{0, 1\}$  based on its content. If the relevance verifier generates “not relevant”, for example, we interpret this as a score of 0. Importantly, to verify the generalizability of our proposed methodology, we keep the implementation of our verifiers fixed for all reasoning tasks and for all datasets. We provide an illustrative example of the verifiers we use in Figure 2. We provide in-context exemplars for the Mathematical Accuracy verifier due to challenges in the LLM’s ability to generate a valid intermediate structured output, a requirement of this verifier’s setup.<sup>2</sup> In addition to the scores from the aforementioned three verifiers, we use the perplexity score of a reasoning step as an additional verifier, in order to encourage

<sup>1</sup>Prompts available in Appendix C

<sup>2</sup>Same exemplars were used for all Math datasets.

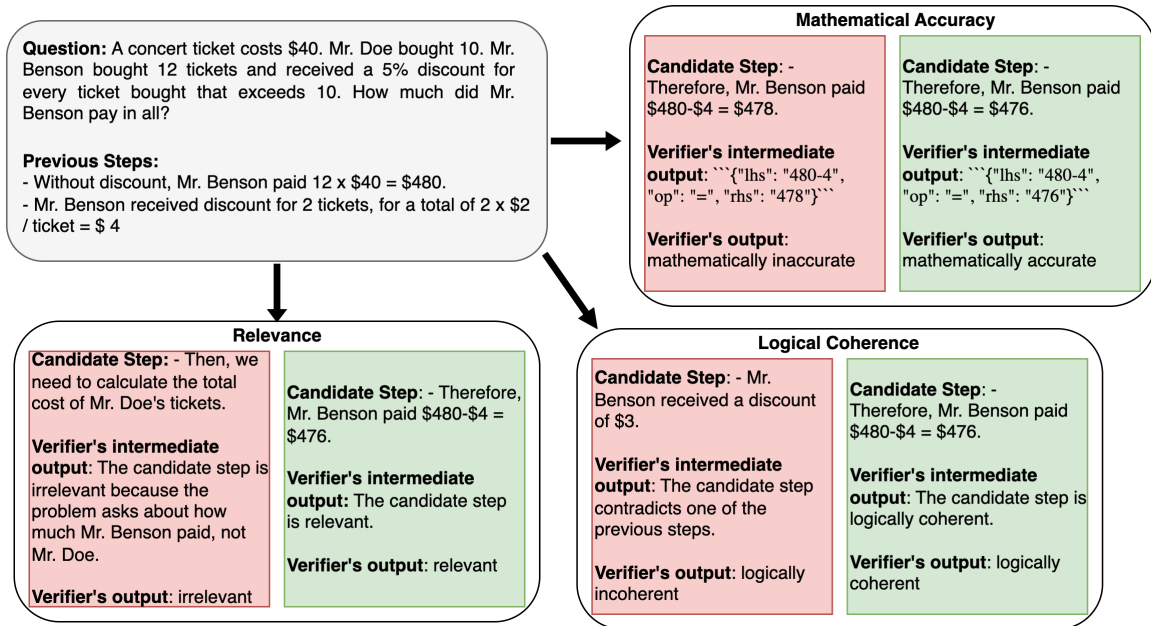


Figure 2: An example of each of our proposed verifiers applied to a given question and previous steps.

the final solution towards text deemed more likely by the LLM. We explain each verifier in greater detail below.

### 3.3.1 Relevance

The first verifier in our proposed framework is the Relevance verifier, where the goal is to constrain a reasoning step to contribute to the construction of a meaningful solution narrative. We provide a *relevant* and an *irrelevant* example in Figure 2. We begin with a question, a sequence of previous steps, and a candidate step. The Relevance verifier assesses the candidate reasoning steps for their relevance to the problem at hand. In the example provided, calculating how much Mr. Doe spent is irrelevant because the problem asks for how much Mr. Benson spent and there is no connection between them.

We acknowledge the inherent subjectivity and nuance associated with determining the relevance of a given reasoning step. However, there are instances where it becomes clear that a reasoning step is distinctly irrelevant, deviating from the coherent solution narrative. For instance, some reasoning steps may veer into speculative or unrelated content, which our Relevance verifier aims to identify.

### 3.3.2 Mathematical Accuracy

The Mathematical Accuracy constraint enforces the need for each reasoning step to contain correct mathematical calculations. We implement this in a similar manner to Tool-based approaches (Schick et al., 2023), working as follows. First, we extract

the mathematical formulas (if present) from a sentence as structured output, as depicted in the *Verifier's intermediate output* field corresponding to the Mathematical Accuracy constraint in Figure 2. For each mathematical calculation present, we extract the left-hand side (*lhs*), the right-hand side (*rhs*), and the operator (*op*). Then, we programmatically execute the extracted formulas (if any) and compare them using the extracted operator.

### 3.3.3 Logical Consistency

A third condition for a logically sound argument we use is for the reasoning steps to not contradict each other (M. and Mckeon, 1941). To this end, we introduce Logical Consistency as our third verifier. This verifier operates over the previous steps and the current candidate step. For example, in Figure 2, the candidate step - *Mr. Benson received a discount of \$3.* contradicts one of the previous steps, as one of the previous steps already established that Mr. Benson received a discount of \$4.

### 3.3.4 Step-wise Perplexity

In addition to the scores resulting from our previously introduced constraint verifiers, we leverage step-wise perplexity as another source of signal, with the goal of favoring lower-perplexity solutions. For each reasoning step  $r_i = [t_1, \dots, t_n]$ , we compute the perplexity over its token constituents. We hypothesize that lower-perplexity reasoning steps are more desirable, as a lower perplexity prompt is correlated with a higher final performance (Gonen et al., 2022). We can interpret a partial reasoning

chain  $R_{1:i} = [r_1, \dots, r_i]$  as (part of) a prompt that will be used to generate  $r_{i+1}$ , making the findings in Gonen et al. (2022) applicable for our purposes.

### 3.4 Constraint Satisfaction Score

Except for the Perplexity verifier, all our proposed verifiers output a binary value, representing whether a given reasoning step satisfies the given constraint or not. For example, if the Relevance verifier gives a score of 1 for a given reasoning step  $r$ , this means that the given step is deemed as relevant. Since the underlying implementation in order to reduce the variance and get a more fine-grained score  $s$ , we use the expected value:  $s = \mathbb{E}(\mathbb{1}_v(r))$ , which we approximate using sampling. Since each verifier is implemented with an LLM, we can sample multiple generations, map each one to a binary value  $\{0, 1\}$ , and then average.

### 3.5 Using Verifiers in Aggregate

#### 3.5.1 Scoring a Reasoning Chain $R$

Given a verifier  $v \in V$ , we extend the concept of a score for a given reasoning step  $r$  to the score for a given (partial or not) reasoning chain  $R$  by aggregating the scores over each of its constituent reasoning steps. Formally, we extend the verifier’s scores to that of a reasoning chain  $R = [r_1, \dots, r_i]$ , where we first obtain a score for each  $r_i$ , resulting in the following score vector:  $[\mathbb{E}(\mathbb{1}_v(r_1)), \dots, \mathbb{E}(\mathbb{1}_v(r_i))]$ , and then aggregate. A low-scoring reasoning step does not necessarily render the entire reasoning chain wrong, but it does increase the likelihood of inaccuracies. To combine these scores, we employ the geometric mean as a milder alternative to the *min* operator in our aggregation process:  $v(R) = GM([\mathbb{E}(\mathbb{1}_v(r_1)), \dots, \mathbb{E}(\mathbb{1}_v(r_i))])$ . We obtain a single score for a given reasoning chain  $R$  and a set of verifiers  $\mathcal{V}$  by aggregating over the scores of each verifier  $v \in \mathcal{V}$  on  $R$ . Our proposed framework allows for the customization of each verifier’s contribution during aggregation. We use a weighted arithmetic mean, as defined below.

$$\mathcal{V}(R) = \frac{\sum_{i=1}^{|\mathcal{V}|} w_i \times v_i(R)}{\sum_{i=1}^{|\mathcal{V}|} w_i}$$

We set  $w = 2$  for perplexity and  $w = 1$  for all the others. We selected  $w = 2$  for perplexity based on preliminary experiments on the train partition of GSM8k and CSQA 2.0. Importantly, we use the same weights for all our experiments.

### 3.5.2 Ensembling Methods using the Verifiers

Ensembling techniques work by aggregating the solution of multiple reasoning chains to obtain a final solution. For example, Self-Consistency (Wang et al., 2022) randomly samples a given number of reasoning chains, and then performs a majority vote on the final answer. Instead of resorting to a majority voting mechanism over randomly sampled reasoning chains, we propose to leverage the scores obtained from our proposed verifiers to do the selection and the weighting.

## 4 Experiments

### 4.1 Experimental Setting

**Models** We use *Falcon*<sup>3</sup> (Almazrouei et al., 2023) as our base LLM, as it was one of the largest and most capable open-source model family freely available at the time of the experiments.<sup>4</sup> We use the same model for both solution generation and solution verification. For solution generation, we use the zero-shot prompt from Kojima et al. (2022). For verification, we use different prompts for each verifier, available in Appendix C.

**Datasets** We perform experiments spanning 4 reasoning tasks: Math, Commonsense, Symbolic, and Other, and 9 datasets: BigBench Date Understanding (bench authors, 2023) (*Other*), CommonsenseQA (Talmor et al., 2019), CommonsenseQA 2.0 (Talmor et al., 2021) and Strategy (Geva et al., 2021) (*Commonsense*), Coinflip and Last Letter Concatenation (Wei et al., 2022) (*Symbolic*), GSM8k (Cobbe et al., 2021), SVAMP (Patel et al., 2021), and AddSub (Kojima et al., 2022) (*Math*). We show an example from each dataset in Appendix B.1.<sup>5</sup>

We use the standard evaluation metrics as previous work, which is Accuracy score computed between the gold answer and the predicted answer.

**Proposed Method Setting** All our proposed verifiers are dataset-agnostic and we use the same prompts for all our experiments. Due to computational constraints, we use the mathematical accuracy verifier only for the math datasets. In an attempt to minimize the impact of the underlying search strategy for the step-by-step solution, we adopt the following approach: we first sample 40

<sup>3</sup>Specifically, we use *Falcon-40B-Instruct*

<sup>4</sup>Open source according to <https://opensource.org/>

<sup>5</sup>For Last Letter Concatenation, we use only 2 words instead of 4.

	Other	Commonsense			Symbolic		Math		
	BigBench Date	CSQA	CSQA 2.0	Strategy	Coinflip	Last Letter (2)	GSM8k	SVAMP	AddSub
Random Chain	55.77±3.06	47.91±1.22	58.75±1.68	56.03±0.99	58.67±2.04	15.68±1.51	29.23±1.58	38.78±1.26	41.04±1.79
Low PPL Chain	63.69±0.00	48.81±0.00	59.10±0.00	<b>60.90±0.00</b>	<b>72.80±0.00</b>	<b>45.00±0.00</b>	40.50±0.00	53.50±0.00	58.48±0.00
Top Chain wrt Verifiers	<b>69.12±0.21</b>	<b>56.79±0.12</b>	<b>62.16±0.22</b>	57.21±0.17	64.02±0.21	41.64±0.44	<b>45.94±0.30</b>	<b>56.36±0.23</b>	<b>62.34±0.25</b>

Table 1: Comparison between two baselines: (1) Random Chains, and (2) Low PPL Chain, and our proposed method, (3) Top Chain wrt Verifiers. In this setting, we record the performance when selecting *one* reasoning chain, according to each method’s selection criteria. We report Accuracy ( $\uparrow$ ).

reasoning chains for each problem and for each dataset, then use the scores resulting from our proposed verifiers to guide our selection process.

**Baselines** We analyze the verifiers’ contributions by comparing the performance of the proposed method against the following baselines: (i) *Random Chains*, where we use the LLM to sample a solution, using the same prompt as Kojima et al. (2022), and (ii) *best-of N* sampling (Adiwardana et al., 2020; Wang et al., 2022), where we sample a total of 40 reasoning chains and select the one with the lowest perplexity. Our motivation for using these two baselines is two-fold. First, we want to allow both the baselines and our proposed method to operate over the same candidate reasoning chains. Secondly, it has been observed in Wang et al. (2022) that Best-of N sampling performs better than greedy decoding, especially for large N.

**Experiments** We conduct the following experiments: (i) Single chain analysis, where we use the resulting scores of each reasoning chain to select a single reasoning chain, (ii) Self-Consistency, where we aggregate multiple reasoning chains, and (iii) Single chain analysis with incomplete chain scoring, where we only score the reasoning chains based on the an initial % (or number) of the reasoning steps.

## 4.2 Single Chain

In this experiment, we assess how the scores generated by our proposed verifiers are correlated with the likelihood of a reasoning chain reaching the correct final answer. For this purpose, we conduct the following experiment: from the 40 sampled reasoning chains, we select the highest-scoring chain based on our proposed verifiers’ scores. We present our results in Table 1, comparing with two baselines: *Random Chains* and *best-of N* sampling. We make the following remarks.

First, over all datasets, our proposed method performs better than selecting a reasoning chain at random, with improvements ranging from 1.18 points (Strategy) to 25.68 points (Last Letter), with

an average improvement of 12.63.

Second, we remark that our proposed method outperforms the best reasoning chain according to perplexity (Low PPL Chains) in over 65% of the cases (6 out of 9 datasets). This means that our proposed verification procedure provides valuable information beyond what is captured by simply selecting the lowest perplexity chains. We note that this trend does not hold true for Symbolic Reasoning, where for both datasets investigated (*Coinflip* and *Last Letter*) the Low PPL Chain is better than the one selected according to our proposed verifier. An exploration over *Coinflip* revealed that steps where the coin has not been flipped received, on average, a lower relevance score, although this information is relevant. We leave the exploration of better verifiers for Symbolic Reasoning to future work. All in all, the average improvement of our proposed method over the reasoning chain with the lowest perplexity is 1.43 points.

## 4.3 Self-Consistency

In this experiment, we explore how well our proposed method leverages ensemble techniques, particularly Self-Consistency (Wang et al., 2022). We start with the same set of 40 reasoning paths and employ different selection strategies to evaluate their effectiveness: (i) we randomly sample from these paths (*Random Chains*), (ii) we select chains with the lowest perplexity from this set (*Low PPL Chains*), and (iii) we choose the reasoning chains with high scores, as determined by the verifiers, from these 40 paths (*Proposed (weighted)*). Unlike the original majority vote approach, we empirically found that weighting each reasoning chain by their verifier scores yields slightly better results. However, it is worth noting that this improvement does not hold when using only perplexity, as shown in (Wang et al., 2022). We show in Figure 3 the behavior of our proposed method. We make two observations: First, our proposed method is able to leverage ensembling techniques, showing consistent performance gains as the number of reasoning chains increases. Second, we remark that our

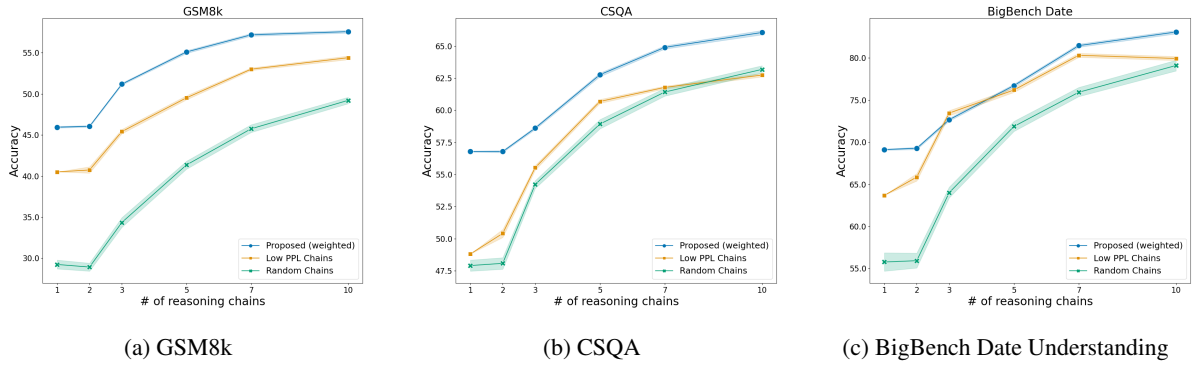


Figure 3: Comparison between our proposed method and two baselines, when using Self-Consistency and between 1-10 reasoning chains. We report Accuracy ( $\uparrow$ ).

proposed method scales better, consistently outperforming the baselines.

We further investigate the performance impact of weighted voting, utilizing scores from our proposed verifiers, against the standard majority-voting approach. Specifically, we apply both voting methods to the *identical set* of reasoning chains, initially selected at random. We found that using the scores of our proposed verifiers to do a weighted voting improves over the majority voting in over 96% of the cases.<sup>6</sup> Due to space constraints, we include the resulting plots in Appendix E.

#### 4.4 Verifying Incomplete Reasoning Chains

In our prior experiments, our proposed verifiers evaluated complete reasoning chains. Now, we explore their effectiveness when applied exclusively to the initial reasoning steps. This experiment provides insights into the potential utility of our proposed method in an “online” setting, where the reasoning step-level evaluation is employed to guide the search for good reasoning chains without fully generating multiple candidate solutions.

We assess the impact of using the verifiers for varying percentages of reasoning steps, denoted as  $X\%$  along the  $X$ -axis of our line plot in Figure 4. We remark that the final performance increases with the  $\%$  of steps verified and that verifying only the first 20% of the steps is sufficient to increase the final performance beyond that of random chains.

Since knowing beforehand the total number of reasoning steps is unrealistic, we also experiment with only verifying a given number of the initial reasoning steps. Due to space limitations we include these results in Appendix G. Additionally, we include in Table 5 the resulting performance when verifying between 0 and *All* reasoning steps

<sup>6</sup>78/81

Verifiers				Math		
P	R	M	C	AddSub	GSM8k	SVAMP
$\times$	$\times$	$\times$	$\times$	41.04 $\pm$ 1.79	29.23 $\pm$ 1.58	38.78 $\pm$ 1.26
$\times$	$\times$	$\times$	$\checkmark$	48.84 $\pm$ 0.85	33.51 $\pm$ 0.69	45.41 $\pm$ 0.56
$\times$	$\times$	$\checkmark$	$\times$	45.76 $\pm$ 1.99	35.59 $\pm$ 1.60	41.68 $\pm$ 1.65
$\times$	$\checkmark$	$\times$	$\times$	49.04 $\pm$ 0.44	33.38 $\pm$ 0.52	46.98 $\pm$ 0.42
$\checkmark$	$\times$	$\times$	$\times$	59.09 $\pm$ 0.59	41.53 $\pm$ 0.49	53.85 $\pm$ 0.45
$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	<b>62.34<math>\pm</math>0.25</b>	<b>45.94<math>\pm</math>0.30</b>	<b>56.36<math>\pm</math>0.23</b>

(a) Ablation over Math datasets.

Verifiers			Other	Commonsense		Symbolic		
P	R	C	BigBench Date	CSQA2.0	CSQA	Strategy	Coinflip	Last Letter
$\times$	$\times$	$\times$	55.77 $\pm$ 3.06	58.75 $\pm$ 1.68	47.91 $\pm$ 1.22	56.03 $\pm$ 0.99	58.67 $\pm$ 2.04	15.68 $\pm$ 1.51
$\times$	$\times$	$\checkmark$	62.70 $\pm$ 0.68	56.77 $\pm$ 0.55	53.00 $\pm$ 0.59	54.31 $\pm$ 0.70	49.17 $\pm$ 0.76	20.63 $\pm$ 0.78
$\times$	$\checkmark$	$\times$	62.57 $\pm$ 0.45	61.03 $\pm$ 0.21	52.62 $\pm$ 0.29	57.31 $\pm$ 0.26	53.63 $\pm$ 0.51	15.79 $\pm$ 0.30
$\checkmark$	$\times$	$\times$	64.06 $\pm$ 0.60	60.27 $\pm$ 0.25	51.35 $\pm$ 0.36	<b>60.03<math>\pm</math>0.30</b>	<b>73.36<math>\pm</math>0.42</b>	<b>41.73<math>\pm</math>0.48</b>
$\checkmark$	$\checkmark$	$\checkmark$	<b>69.12<math>\pm</math>0.21</b>	<b>62.16<math>\pm</math>0.22</b>	<b>56.79<math>\pm</math>0.12</b>	57.21 $\pm$ 0.17	64.02 $\pm$ 0.21	41.64 $\pm$ 0.44

(b) Ablation over Non-Math datasets.

Table 2: Ablation study on the effect of each verifier on the downstream tasks when selecting a *single* reasoning chain. We differentiate between math and non-math datasets.

over all datasets. In 7/9 cases, the performance increases even when verifying only the first step. When verifying the first two steps, the final performance increases in all the cases.

#### 4.5 Contributions of each Verifier

In this experiment, we assess the contribution to the final performance of each of our verifiers: (1) Low Step Perplexity, (2) **R**elevance, (3) **M**athematical Accuracy (if applicable), (4) Logical **C**onsistency. First, we observe that each individual verifier is meaningfully contributing towards the final solution. For example, for the math datasets (Table 2a), employing *any* verifier improves the final performance, with improvements ranging from 2.90% to 21.30%. All in all, using as little as a single verifier improves the final performance in over 89% of the cases.<sup>7</sup> Secondly, we remark that combining all the verifiers gives further improvements, beyond those obtained by using a single verifier,

<sup>7</sup>35/39

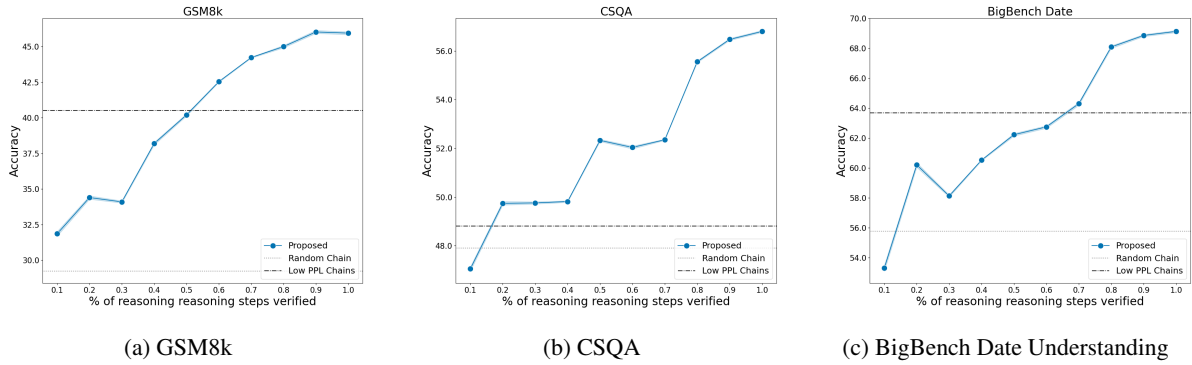


Figure 4: Verifying only the first X% steps of a given reasoning chain. (†).

suggesting that each verifier is adding meaningful and non-overlapping information. We note that there is a notable exception to this trend, where for the Symbolic Reasoning tasks (and for the Strategy dataset), a distinct combination of verifiers (i.e. only Perplexity) attains a better score than using all the verifiers. We provide results covering a wider range of verifier combinations in Appendix F.

#### 4.6 Human Evaluation

While the proposed verifiers meaningfully contribute to the final performance on the reasoning downstream tasks, we perform a human evaluation study to assess: (1) how well they correlate with human judgment and (2) how reliably concepts such as logical consistency or relevance can be evaluated by humans. We include instructions and inter-annotator agreement scores in Appendix H.

We compute Pearson correlation scores (in range  $[-1, 1]$ ) between human assessments and the scores proposed by our verifiers and, additionally, also between GPT-4o and humans to evaluate how the proposed method scales with a stronger underlying model. We plot average correlations up to step K in Figure 5. We note the following.

First, we note that each verifier exhibits an overall significant ( $p$ -value  $< 0.0001$ ) and positive correlation between human judgement and performance. When the correlation is negative, it is better to not use it (e.g., Coinflip has scores 69 (without) vs 64 (with) the negatively correlated logical consistency verifier, see 4b). We also show in Appendix I how even small positive correlations statistically differentiate better outcomes on average. Second, we observe (and show in Appendix H) a large variance in the inter-annotator agreement score, which we hypothesize comes from different reasoning styles between humans and noise from a weak underlying model (Falcon). We also explore different agreement metrics, aiming to capture some of the

Cohen’s Kappa Limitations (e.g., Cohen’s Kappa Paradox (Zec et al., 2017)). Lastly, we found that less than 2% of the errors marked by the annotators are not captured by one of the principles explored.

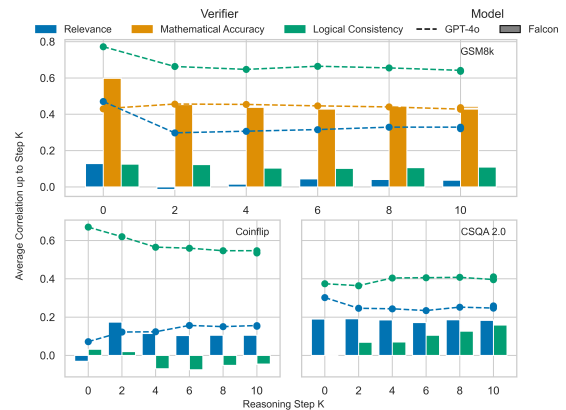


Figure 5: Correlation across datasets and steps

## 5 Conclusion

We explore a general-purpose verification procedure consisting of task- and dataset-agnostic verifiers at the reasoning step-level inspired by fundamental principles of sound reasoning: Relevance, Mathematical Accuracy and Logical consistency. On top of these reasoning principles, we leverage the perplexity of the reasoning step to steer the LLM towards high-quality solutions.

Across four distinct reasoning tasks, spanning nine datasets, we show that using the proposed verifiers to score the reasoning chains leads to notable performance improvements when compared to randomly sampled reasoning chains. Most notably, our proposed approach outperforms the lowest perplexity reasoning chain in over 6 out of the 9 datasets we tested. This indicates that the proposed verifiers provide additional valuable information beyond what is captured by the perplexity measure of the reasoning chain. We leave the exploration of better and more efficient verifiers to future work.



## 598 Limitations

599 While our proposed framework is flexible and ad-  
600 mits different implementations for the verifiers, in  
601 this work we implemented each verifier with a  
602 prompt-based LLM approach. This type of im-  
603 plementation can increase the energy consumption  
604 of the deployed system, leading to a performance-  
605 energy-consumption trade-off. Secondly, employ-  
606 ing step-by-step verifiers increases the computa-  
607 tional time needed by the system to produce an out-  
608 put. This trade-off must be analyzed on a case-by-  
609 case basis and compared to other alternatives, such  
610 as self-consistency. Different from self-consistency,  
611 our proposed approach aims to improve the correct-  
612 ness of each step in a step-by-step solution.

613 While our evaluation spanned multiple reason-  
614 ing tasks and datasets, it was limited to tasks in the  
615 English language only. We leave the evaluation on  
616 more challenging datasets (e.g., MATH (Hendrycks  
617 et al., 2021)) or datasets with contradictory infor-  
618 mation (Chen et al., 2022; Kazemi et al., 2023) to  
619 future work.

620 During our human evaluation study, we observed  
621 low, yet significant, positive correlations. To vali-  
622 date the observed performance improvements on  
623 downstream tasks, we conducted additional experi-  
624 ments on synthetic data, which confirmed that these  
625 improvements are indeed expected.

626 Lastly, we only used Falcon-40B-Instruct in our  
627 experiments, as it was one of the largest and most  
628 capable open-source<sup>8</sup> model at the time of the ex-  
629 periments. We utilized GPT-4o and human corre-  
630 lations to provide a glimpse into how verifiers im-  
631 plemented with a more powerful underlying model  
632 might perform in comparison. We leave the explo-  
633 ration of other models to future work.

## 634 Ethics Statement

635 Our proposed approach utilizes large language  
636 models, which are known to be biased and to hallu-  
637 cinate. In this work, we do not pre-train nor fine-  
638 tune any large-scale models. Instead, we use al-  
639 ready pre-trained open-source models and prompt-  
640 ing.

## 641 References

642 Daniel Adiwardana, Minh-Thang Luong, David R. So,  
643 Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang,  
644 Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu,

and Quoc V. Le. 2020. [Towards a human-like open-domain chatbot.](#) 645 646

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-  
shamsi, Alessandro Cappelli, Ruxandra-Aimée Co-  
jocararu, Daniel Hesslow, Julien Launay, Quentin  
Malartic, Daniele Mazzotta, Badreddine Noune, Bap-  
tiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models.](#) *ArXiv*,  
abs/2311.16867. 647 648 649 650 651 652 653

Yuntao Bai, Saurav Kadavath, Sandipan Kundu,  
Amanda Askell, Jackson Kernion, Andy Jones, Anna  
Chen, Anna Goldie, Azalia Mirhoseini, Cameron  
McKinnon, Carol Chen, Catherine Olsson, Christo-  
pher Olah, Danny Hernandez, Dawn Drain, Deep  
Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez,  
Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua  
Landau, Kamal Ndousse, Kamile Lukosuite, Liane  
Lovitt, Michael Sellitto, Nelson Elhage, Nicholas  
Schiefer, Noemi Mercado, Nova DasSarma, Robert  
Lasenby, Robin Larson, Sam Ringer, Scott John-  
ston, Shauna Kravec, Sheer El Showk, Stanislav Fort,  
Tamera Lanham, Timothy Telleen-Lawton, Tom Con-  
erly, Tom Henighan, Tristan Hume, Samuel R. Bow-  
man, Zac Hatfield-Dodds, Ben Mann, Dario Amodei,  
Nicholas Joseph, Sam McCandlish, Tom Brown, and  
Jared Kaplan. 2022. [Constitutional ai: Harmlessness from ai feedback.](#) 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671

BIG bench authors. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.](#) *Transactions on Machine Learning Research.* 672 673 674 675

Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-  
stenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz  
Lehmann, Michal Podstawski, Hubert Niewiadom-  
ski, Piotr Nyczyk, and Torsten Hoefler. 2023. [Graph of thoughts: Solving elaborate problems with large language models.](#) 676 677 678 679 680 681

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, T. J. Henighan, Rewon Child,  
Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens  
Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
Mateusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners.](#) *ArXiv*,  
abs/2005.14165. 682 683 684 685 686 687 688 689 690 691 692 693

Sébastien Bubeck, Varun Chandrasekaran, Ronen El-  
dan, John A. Gehrke, Eric Horvitz, Ece Kamar, Peter  
Lee, Yin Tat Lee, Yuan-Fang Li, Scott M. Lundberg,  
Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro,  
and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4.](#) *ArXiv*,  
abs/2303.12712. 694 695 696 697 698 699 700

Hung-Ting Chen, Michael J.Q. Zhang, and Eunsol Choi.  
2022. [Rich knowledge sources bring complex knowl-](#) 701 702

<sup>8</sup>Open source according to <https://opensource.org/>

703	<a href="#">edge conflicts: Recalibrating models to reflect conflicting evidence.</a> In <i>Conference on Empirical Methods in Natural Language Processing</i> .	
704		
705		
706	Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023. <a href="#">Teaching large language models to self-debug.</a> <i>ArXiv</i> , abs/2304.05128.	
707		
708		
709	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pilla, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. <a href="#">Palm: Scaling language modeling with pathways.</a> <i>J. Mach. Learn. Res.</i> , 24:240:1–240:113.	
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. <a href="#">Training verifiers to solve math word problems.</a> <i>ArXiv</i> , abs/2110.14168.	
734		
735		
736		
737		
738		
739	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. <a href="#">Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies.</a> <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	
740		
741		
742		
743		
744		
745	Hila Gonen, Srini Iyer, Terra Blevins, Noah A. Smith, and Luke Zettlemoyer. 2022. <a href="#">Demystifying prompts in language models via perplexity estimation.</a> <i>ArXiv</i> , abs/2212.04037.	
746		
747		
748		
749	Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2023. <a href="#">Think before you speak: Training language models with pause tokens.</a> <i>ArXiv</i> , abs/2310.02226.	
750		
751		
752		
753		
754	Kilem L. Gwet. 2014. <a href="#">Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters.</a>	
755		
756		
757	Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. <a href="#">Reasoning with language model is planning with world model.</a> <i>ArXiv</i> , abs/2305.14992.	
758		
759		
760		
	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. <a href="#">Measuring mathematical problem solving with the math dataset.</a> <i>NeurIPS</i> .	761 762 763 764
	Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023. <a href="#">Large language models cannot self-correct reasoning yet.</a> <i>ArXiv</i> , abs/2310.01798.	765 766 767 768 769
	Mehran Kazemi, Quan Yuan, Deepti Bhatia, Najoung Kim, Xin Xu, Vaiva Imbrasaite, and Deepak Ramachandran. 2023. <a href="#">Boardgameqa: A dataset for natural language reasoning with contradictory information.</a> <i>ArXiv</i> , abs/2306.07934.	770 771 772 773 774
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. <a href="#">Large language models are zero-shot reasoners.</a> <i>ArXiv</i> , abs/2205.11916.	775 776 777 778
	Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023. <a href="#">Making language models better reasoners with step-aware verifier.</a> In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.	779 780 781 782 783 784 785
	Jieyi Long. 2023. <a href="#">Large language model guided tree-of-thought.</a> <i>ArXiv</i> , abs/2305.08291.	786 787
	E. A. M. and Richard Mckeon. 1941. <a href="#">The basic works of aristotle.</a>	788 789
	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder, Shashank Gupta, Amir Yazdanbakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Iterative refinement with self-feedback.</a> <i>ArXiv</i> , abs/2303.17651.	790 791 792 793 794 795 796
	Niels Mündler, Jingxuan He, Slobodan Jenko, and Martin T. Vechev. 2023. <a href="#">Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation.</a> <i>ArXiv</i> , abs/2305.15852.	797 798 799 800
	Maxwell Nye, Anders Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. 2021. <a href="#">Show your work: Scratchpads for intermediate computation with language models.</a> <i>ArXiv</i> , abs/2112.00114.	801 802 803 804 805 806 807
	OpenAI. 2023. <a href="#">Gpt-4 technical report.</a> <i>ArXiv</i> , abs/2303.08774.	808 809
	Liangming Pan, Michael Stephen Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2023. <a href="#">Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies.</a> <i>ArXiv</i> , abs/2308.03188.	810 811 812 813 814

815	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	Geoffrey Irving, and Irina Higgins. 2022. <a href="#">Solving math word problems with process- and outcome-based feedback</a> . <i>ArXiv</i> , abs/2211.14275.	871
816	2021. <a href="#">Are NLP models really able to solve simple math word problems?</a> In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2080–2094, Online. Association for Computational Linguistics.		872
817			873
818		Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. 2022. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . <i>ArXiv</i> , abs/2203.11171.	874
819			875
820			876
821			877
822	William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Ouyang Long, Jonathan Ward, and Jan Leike. 2022. <a href="#">Self-critiquing models for assisting human evaluators</a> . <i>ArXiv</i> , abs/2206.05802.	Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023. <a href="#">Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents</a> . <i>ArXiv</i> , abs/2302.01560.	878
823			879
824			880
825			881
826	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. <a href="#">Toolformer: Language models can teach themselves to use tools</a> . <i>ArXiv</i> , abs/2302.04761.		882
827		Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. <a href="#">Chain of thought prompting elicits reasoning in large language models</a> . <i>ArXiv</i> , abs/2201.11903.	883
828			884
829			885
830			886
831	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Huai hsin Chi, Nathanael Scharli, and Denny Zhou. 2023. <a href="#">Large language models can be easily distracted by irrelevant context</a> . In <i>International Conference on Machine Learning</i> .		887
832		Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, MingSung Kan, Junxian He, and Qizhe Xie. 2023. <a href="#">Self-evaluation guided beam search for reasoning</a> .	888
833			889
834			890
835			
836	Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. <a href="#">Reflexion: Language agents with verbal reinforcement learning</a> .	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. <a href="#">Tree of thoughts: Deliberate problem solving with large language models</a> . <i>ArXiv</i> , abs/2305.10601.	891
837			892
838			893
839			894
840			895
841	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. <a href="#">CommonsenseQA: A question answering challenge targeting commonsense knowledge</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. <a href="#">React: Synergizing reasoning and acting in language models</a> . <i>ArXiv</i> , abs/2210.03629.	896
842			897
843			898
844			899
845		Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. <a href="#">Do large language models know what they don't know?</a> In <i>Annual Meeting of the Association for Computational Linguistics</i> .	900
846			901
847			902
848			903
849	Alon Talmor, Ori Yoran, Ronan Le Bras, Chandrasekhar Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2021. <a href="#">Commonsenseqa 2.0: Exposing the limits of ai through gamification</a> . <i>ArXiv</i> , abs/2201.05320.	Slavica Zec, Nicola Soriani, Rosanna Irene Comoretto, and Ileana Baldi. 2017. <a href="#">High agreement and high prevalence: The paradox of cohen's kappa</a> . <i>The Open Nursing Journal</i> , 11:211 – 218.	904
850			905
851			906
852			907
853			908
854	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open and efficient foundation language models</a> . <i>ArXiv</i> , abs/2302.13971.		
855			
856			
857			
858			
859			
860			
861	Miles Turpin, Julian Michael, Ethan Perez, and Sam Bowman. 2023. <a href="#">Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting</a> . <i>ArXiv</i> , abs/2305.04388.		
862			
863			
864			
865	Gladys Tyen, Hassan Mansoor, Peter Chen, Tony Mak, and Victor Carbune. 2023. <a href="#">Llms cannot find reasoning errors, but can correct them!</a> <i>ArXiv</i> , abs/2311.08516.		
866			
867			
868			
869	Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, L. Wang, Antonia Creswell,		
870			

909	<b>A Frequently Asked Questions</b>	
910	<b>A.1 What performance improvements are expected given certain correlation levels?</b>	
911		
912	We argue that the low correlations are not problematic and that a correlation of 1.0 is (currently) unrealistic, as it would imply a verifier as good as humans. To further support this assertion and to add details on the level of improvements that can be expected from correlations like 0.1, we ran experiments on artificially generated data and show the results in Table 10. This experiment further demonstrates that the improvements we have seen with correlation scores of 0.1 are normal. Intuitively, if the scores of the verifiers are not correlated with humans for a given question, the worst it can do is to select a random chain. However, if the verifiers are correlated with humans for a given reasoning chain, then they will select those reasoning chains that humans agree are better.	
913		
914		
915		
916		
917		
918		
919		
920		
921		
922		
923		
924		
925		
926		
927		
928	We supplement the correlation values of Falcon with that of GPT-4o, a more powerful model, to highlight how the proposed method scales with a stronger underlying model. Overall, GPT-4o obtains better correlation scores with human judgments.	
929		
930		
931		
932		
933		
934	<b>A.2 The variance in Annotator’s agreement scores</b>	
935		
936	We hypothesize that the large variance in the inter-annotator agreement score comes from different reasoning styles between different humans. Nevertheless, the overall agreements presented in Table 6 highlight between Moderate Agreement and Substantial Agreement. We show all the inter-annotator agreements over each of the four attributes in Appendix H. We also elaborate on how Cohen’s Kappa might be too harsh, a phenomenon known as Cohen’s Kappa Paradox (Zec et al., 2017). For example, the Cohen’s Kappa scores between the following two annotations $a_1 = [1, 1, 1, 1, 1, 1]$ , $a_2 = [1, 1, 1, 1, 1, 0]$ is 0, even though they only disagree on one instance.	
937		
938		
939		
940		
941		
942		
943		
944		
945		
946		
947		
948		
949		
950	To this end, we included additional agreement scores: Gwet’s AC1 (Gwet, 2014) and naive agreement. Both showed that the annotators agree more than initially revealed by Cohen’s Kappa.	
951		
952		
953		
954	<b>A.3 Costs of the proposed method</b>	
955	Our focus in this work has been to explore the extent to which LLMs are capable of error detection and error correction. Especially since there is a	
956		
957		
	divergence in findings across prior work regarding the LLMs ability to detect and correct its errors (Tyen et al., 2023; Huang et al., 2023; Madaan et al., 2023, <i>inter alia</i> ) even if it meant temporarily sacrificing efficiency. We leave the exploration of more efficient methods to future work.	958 959 960 961 962 963
	<b>A.4 Subjectivity of the verifiers</b>	964
	The relevance and logical consistency verifiers may appear subjective, evidenced by the varying responses from human annotators. However, our focus has not been on addressing nuanced scenarios where human judgment itself might differ. Instead, our intention has been to capture and rectify the more straightforward errors, those instances where even human annotators unanimously agree.	965 966 967 968 969 970 971 972
	<b>B Experimental Settings</b>	973
	<b>B.1 Datasets</b>	974
	We include in Table 3 an input/output example for each dataset used. We also experimented with the Object Tracking problem from BigBench, but we removed it because for all experimental settings, encompassing both baseline and proposed methods, the performance consistently fell below the chance-level threshold. For Last Letter Concatenation, we use only 2 words instead of 4, as we empirically observed that Falcon-40B is not able to tackle the problem when there are four words. We use a temperature of 0.7 for all our experiments.	975 976 977 978 979 980 981 982 983 984 985
	<b>B.2 Hardware</b>	986
	We ran the experiments on machines with 8 A10G 24GB GPUs (AWS g5.48xlarge). In total, we used approximately 10 weeks worth of g5.48xlarge time.	987 988 989 990
	<b>C Verifiers</b>	991
	We include the prompts we used for each verifier.	992
	<b>C.1 Relevance Prompt</b>	993
	You are a helpful assistant that is good at evaluating reasoning chains in order to solve logic problems.	994 995 996
	You are given a logic problem and a draft solution with numbered steps that we need to complete. Evaluate the draft solution by determining whether it adds relevant information that helps to solve the problem. If it is relevant answer by ‘yes, the solution is relevant’, otherwise	997 998 999 1000 1001 1002 1003

<b>Reasoning Task</b>	<b>Dataset Name</b>	<b>Input</b>	<b>Expected Output</b>
Other	BigBench Date Understanding	Yesterday was April 30, 2021. What is the date today in MM/DD/YYYY? Choices: (A) 05/01/2021, (B) 02/23/2021, (C) 03/11/2021, (D) 05/09/2021, (E) 06/12/2021, (F) 04/29/2021	(A)
Commonsense	CSQA 2.0 Strategy	a pupil can be either a student or part of an eye Is it common to see frost during some college commencements?	yes yes
Symbolic	Coinflip	A coin is heads up. Whitney flips the coin. Erika does not flip the coin. Tj does not flip the coin. Benito flips the coin. Is the coin still heads up? Note that "flip" here means "reverse"	yes
	Last Letter (2)	Take the last letters of each words in "Whitney Benito" and concatenate them.	yo
Math	GSM8k	Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May. How many clips did Natalia sell altogether in April and May?	72
	SVAMP	Bryan took a look at his books as well. If he has 34 books distributed equally in 2 bookshelves. How many books are there in each bookshelf?	17
	AddSub	Joan found 70 seashells on the beach. she gave Sam some of her seashells. She has 27 seashell. How many seashells did she give to Sam ?	43

Table 3: An example of input and expected output for each of the datasets we experiment with.

1004	say 'no' and explain which steps failed.		1053
1005	### Problem: {problem statement}	Last Step:	1054
1006	### Draft solution: {previous steps}	{current step}	1055
1007	### Draft step: {current step}		1056
1008	### Your evaluation of the draft step:	Instruction:	1057
1009	<b>C.2 Mathematical Accuracy Prompt</b>	Given the information present in the	1058
1010	For brevity, we have included a single example	Last Step and in the Previous Steps,	1059
1011	from the set of 20 in-context examples. The in-	please check if the conclusion present	1060
1012	context examples are a mix of inputs with math-	in the Last Step is contradicting any	1061
1013	ematical calculations and without. Furthermore,	information from the Previous Steps.	1062
1014	within the set of examples involving mathematical		1063
1015	calculations, both incorrect and correct calculations	Feedback:	1064
1016	are included for comprehensive coverage.	Based on the Last Step, which is	1065
1017		"{current step}", and on the Previous	1066
1018	Instruction:	Steps, we can conclude that the Last	1067
1019	The task is to extract the mathematical	Step is""	1068
1020	calculations appearing below and return		
1021	the result in JSON format. Please do	<b>D Self-Consistency</b>	1069
1022	not perform any additional calculations	We include here the plots for all the datasets for	1070
1023	and do not introduce any number or	the Self-Consistency experiment performed in Sec-	1071
1024	numerical expression that does not	tion 4.3	1072
1025	appear in the original input text. If there		
1026	is no explicit calculation performed, do	<b>E Self-Consistency: Weighted Voting vs</b>	1073
1027	not return anything.	<b>Majority Voting</b>	1074
1028		We include in Figure 7 the results over all datasets	1075
1029	Input:	for the (1) Weighted Voting and (2) Majority Voting	1076
1030	Therefore, he has $\$87-\$32=\ll 87-$	over the same reasoning chains.	1077
1031	$32=40\rangle\$40$ left		
1032	Output:	<b>F Ablation</b>	1078
1033	““json	We include a more comprehensive ablation in Ta-	1079
1034	{[{"lhs": "87-32", "op": "=", "rhs":	ble 4, where we ablate over more combinations of	1080
1035	"40"}]}	the verifiers.	1081
1036	““		
1037	<..> Input:	<b>G Verifying Incomplete Reasoning</b>	1082
1038	{input}	<b>Chains</b>	1083
1039		To complement Figure 4, where we verified a <i>per-</i>	1084
1040	Output:	<i>centage</i> of the total number of reasoning steps for	1085
1041	““json	a given reasoning chain, we also include Figure 8,	1086
1042		where we verify a fixed number of reasoning steps.	1087
1043	<b>C.3 Logical Consistency Prompt</b>	The motivation behind this experiment is that the	1088
1044	""You are a smart, critical, and logical	number of reasoning steps might not be known	1089
1045	teacher assistant. You are critically	beforehand.	1090
1046	reading a student's answer line by	To supplement the analysis performed in Sec-	1091
1047	line and verifying each line for any	tion 4.4, we include in Table 5 the performance	1092
1048	contradictions in the student's argument.	obtained by our proposed method when verifying a	1093
1049	More information below.	varying number of steps, from 0 (no verification)	1094
1050		to All (verify all reasoning steps).	1095
1051	Previous Steps:		
1052	{previous steps}		

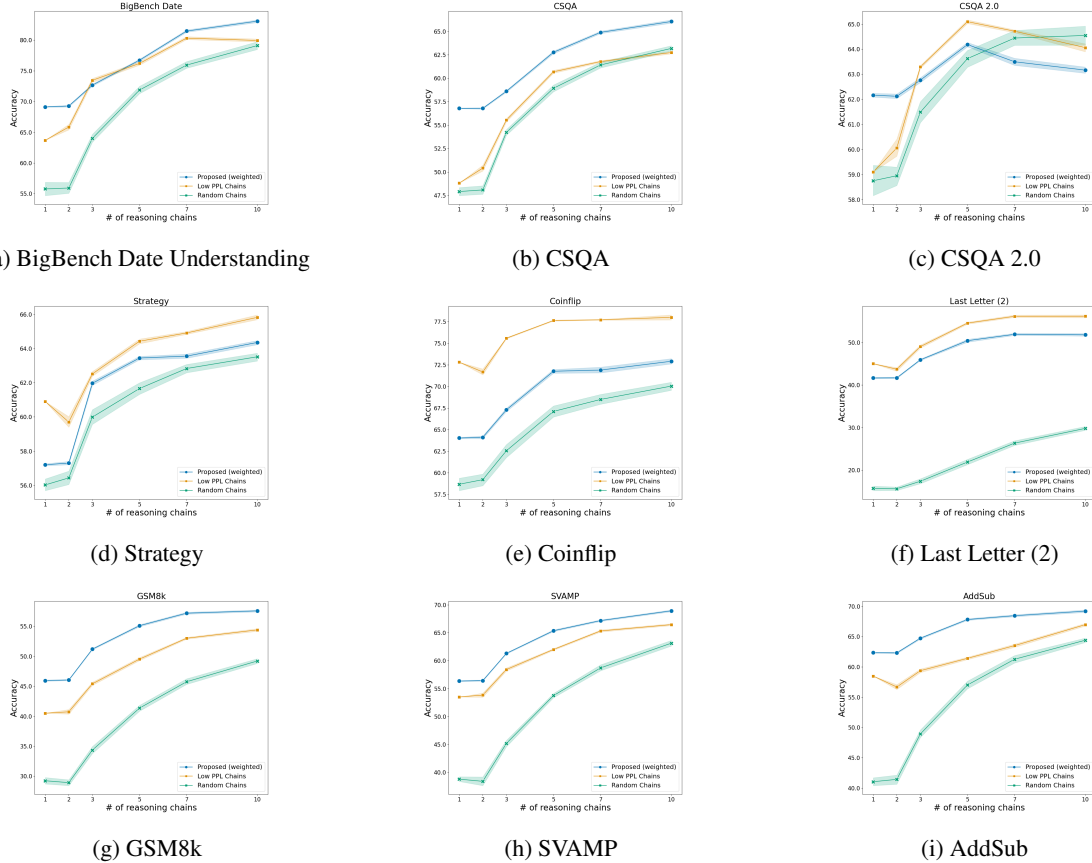


Figure 6: Self-consistency, sampling between 1-10 reasoning paths (↑)

## H Human Evaluation

### H.1 Inter-Annotator Agreement

We include the inter-annotator agreement across all 4 measured attributes in Figures 9, 10, 11, 12. We remark the large variance in the agreement, even for principles that are less subjective (e.g. *Mathematical Accuracy*)

We also include in Table 6 the overall agreement by attribute.

### H.2 Other Agreement Scores

We further investigated the human annotation data and found many instances where the agreement score, as given by Cohen’s Kappa, was too harsh. For example, the Cohen’s Kappa scores between the following two annotations  $a_1 = [1, 1, 1, 1, 1, 1]$  and  $a_2 = [1, 1, 1, 1, 1, 1]$  is *nan*. Between  $a_3 = [1, 1, 1, 1, 1, 1]$ ,  $a_4 = [1, 1, 1, 1, 1, 0]$  is 0, even though they only disagree on one instance. Lastly, for  $a_5 = [1, 1, 1, 0, 1, 1, 1, 1, 0, 1]$  and  $a_6 = [1, 1, 1, 1, 1, 0, 1, 1, 1, 1]$  the Cohen’s Kappa score is  $-0.154$ . We remark that in all three cases above, judging only from the annotations, the

annotators tend to agree, but this is not reflected in the Cohen’s Kappa score, a phenomenon called Cohen’s Kappa Paradox (Zec et al., 2017). We note that Krippendorff’s alpha only fixes the first example.

Therefore, we include in Table 7 the agreement score, as calculated using Gwet’s AC1 (Gwet, 2014). We remark that the agreement scores are higher than initially revealed by Cohen’s Kappa scores.

Lastly, we include in Table 8 the agreement computed as the percentage of time the annotators give the same label, without accounting for agreement by chance.

### H.3 Correlations

We include the Pearson correlations between each of the 4 measured attributes (Relevance, Logical Consistency, Mathematical Accuracy, and Overall Correctness) in Figure 13. We remark that there is a large variance in the correlation scores.

We also include in Figure 14 the correlations between the verifiers and the human assessments. Additionally, we include the correlations between

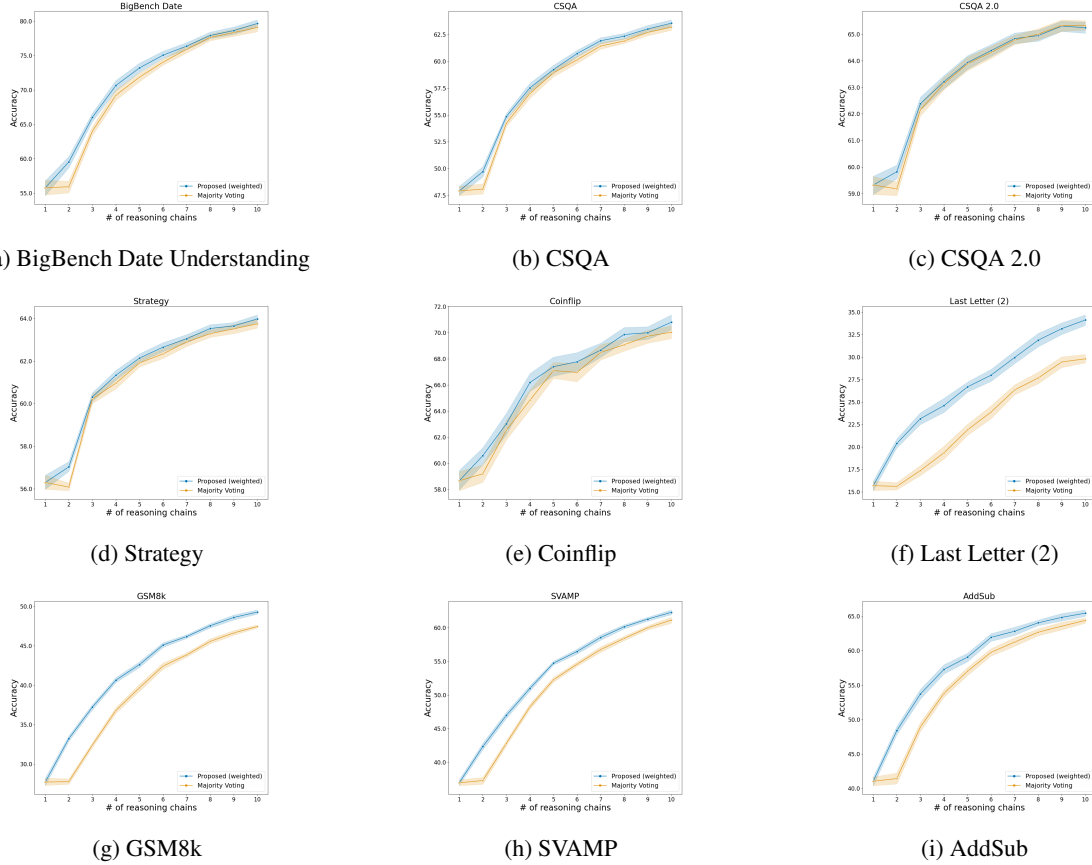


Figure 7: Self-consistency on the *same* reasoning chains, comparing between weighting the final answer using the scores from our proposed verifiers or taking the majority voting ( $\uparrow$ ). Overall, using the scores of our proposed verifiers to perform weighted voting consistently improves the final performance.

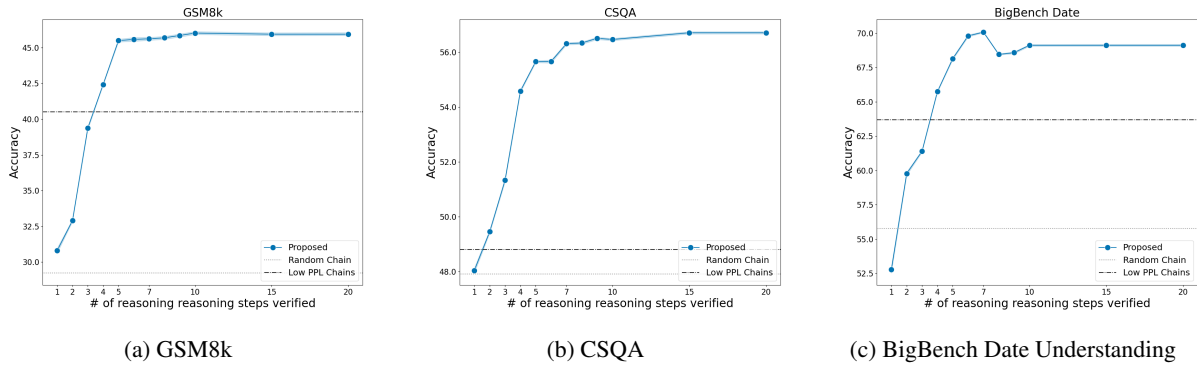


Figure 8: Verifying only the first X steps of a chain.

1141 the human assessment of the overall correctness of  
 1142 a given reasoning step and the aggregated score,  
 1143 with and without perplexity.

#### 1144 H.4 Human Annotator Instructions

1145 We provide an overview of the instructions given  
 1146 to the human annotators.

1147 **(1) Overall Correctness:** If there is any reason-  
 1148 ing issue with this step, answer n. If you can-

not evaluate the step because you lack expertise or  
 there are some other issues with the step, answer a.  
 Please let us know about the reason in the notes col-  
 umn. If there is nothing wrong with the reasoning  
 step, answer y.

**(2) Mathematical Accuracy:** Are all arith-  
 metic calculations in this reasoning step correct?  
 This question is strictly about arithmetic calcula-  
 tions and not about how the calculation is used to

1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157



Verifiers				Math		
Perplexity	Relevance	Math Accuracy	Consistency	AddSub	GSM8k	SVAMP
✗	✗	✗	✗	41.04±1.79	29.23±1.58	38.78±1.26
✗	✗	✗	✓	48.84±0.85	33.51±0.69	45.41±0.56
✗	✗	✓	✗	45.76±1.99	35.59±1.60	41.68±1.65
✗	✓	✗	✗	49.04±0.44	33.38±0.52	46.98±0.42
✗	✓	✓	✓	53.21±0.40	43.37±0.25	49.83±0.18
✓	✗	✗	✗	59.09±0.59	41.53±0.49	53.85±0.45
✓	✗	✓	✗	<u>60.58±0.47</u>	45.66±0.31	55.80±0.37
✓	✗	✓	✓	59.57±0.21	<b>48.02±0.24</b>	55.35±0.24
✓	✓	✗	✗	57.47±0.33	41.87±0.29	54.83±0.15
✓	✓	✓	✗	59.82±0.22	44.47±0.22	<b>56.44±0.24</b>
✓	✓	✓	✓	<b>62.34±0.25</b>	<u>45.94±0.30</u>	<u>56.36±0.23</u>

(a) Ablation Single Chain math results, Accuracy, Higher is better ↑

Verifiers			Other	Commonsense			Symbolic	
Perplexity	Relevance	Consistency	BigBench Date	CSQA 2.0	CSQA	Strategy	Coinflip	Last Letter (2)
✗	✗	✗	55.77±3.06	58.75±1.68	47.91±1.22	56.03±0.99	58.67±2.04	15.68±1.51
✗	✗	✓	62.70±0.68	56.77±0.55	53.00±0.59	54.31±0.70	49.17±0.76	20.63±0.78
✗	✓	✗	62.57±0.45	61.03±0.21	52.62±0.29	57.31±0.26	53.63±0.51	15.79±0.30
✗	✓	✓	63.50±0.26	60.47±0.12	<u>54.74±0.25</u>	55.98±0.24	52.68±0.37	16.50±0.24
✓	✗	✗	64.06±0.60	60.27±0.25	51.35±0.36	<b>60.03±0.30</b>	<b>73.36±0.42</b>	<u>41.73±0.48</u>
✓	✗	✓	64.72±0.31	58.74±0.12	52.99±0.24	58.54±0.13	64.89±0.38	<b>46.03±0.34</b>
✓	✓	✗	<b>71.04±0.22</b>	<b>62.21±0.19</b>	54.28±0.14	<u>58.98±0.14</u>	<u>69.81±0.20</u>	39.53±0.21
✓	✓	✓	<u>69.12±0.21</u>	<u>62.16±0.22</u>	<b>56.79±0.12</b>	57.21±0.17	64.02±0.21	41.64±0.44

(b) Ablation Single Chain non math results, Accuracy, Higher is better ↑

Table 4: Ablation over the types of verifiers used. Overall, all verifiers are meaningfully contributing towards the final solution.

1158 progress toward the solution.

1159 **(3) Logical Consistency:** Is this reasoning step  
1160 logically consistent, in itself, and with previous  
1161 steps? Answer y if the step is logically consistent  
1162 within itself, and with all previous steps, including  
1163 the prompt (problem statement). A step is logically  
1164 consistent when it uses available information in a  
1165 way that is logically correct. In most cases this  
1166 means that the conclusions that are reached in this  
1167 step follow logically from assumptions made. It  
1168 can also mean that the step does not contradict  
1169 information provided in previous steps (or the same  
1170 step).

1171 **(4) Relevance:** Does this reasoning step add in-  
1172 formation that is relevant for solving the problem?  
1173 Information is relevant when it is useful for solv-  
1174 ing the problem (e.g. it states helping assumptions,  
1175 or it reaches a conclusion that answers the prob-  
1176 lem or get you closer to an answer). Information  
1177 can be added by re-stating information from the  
1178 prompt, by reaching a conclusion, or by introduc-  
1179 ing completely new information – any of these can

be relevant or irrelevant.

## 1180 H.5 Ratios of Steps Annotated as Incorrect 1181

1182 We show in Table 9 the ratio of steps each annotator  
1183 deemed as incorrect. 1183

## 1184 I Expected Performance Given 1185 1185 Correlation Levels

1186 We showed in Sections 4.2 and 4.3 that employ- 1186  
1187 ing the proposed verifiers leads to performance 1187  
1188 improvements. Then, we analyzed in Section 4.6 1188  
1189 the correlations between the proposed verifiers and 1189  
1190 human judgments, observing significantly positive 1190  
1191 (but low) correlations. In this section, we further 1191  
1192 analyze what improvements can we expect for a given 1192  
1193 level of correlations. To this end, we conducted ad- 1193  
1194 ditional experiments using artificially generated 1194  
1195 data, allowing us precise control over the correla- 1195  
1196 tion values between the verifiers and human judg- 1196  
1197 ments. By randomly sampling scores to simulate 1197  
1198 the verifiers’ and human annotators’ judgments, we 1198  
1199 manipulated the data to induce positive correlations 1199

# of Steps Verified	Other	Commonsense			Symbolic		Math		
	BigBench Date	CSQA 2.0	CSQA	Strategy	Coinflip	Last Letter (2)	GSM8k	SVAMP	AddSub
0	55.77±3.06	58.75±1.68	47.91 ± 1.22	56.03±0.99	58.67±2.04	15.68±1.51	29.23±1.58	38.78±1.26	41.04±1.79
1	52.77±0.66	59.73±0.26	48.04 ± 0.26	<b>58.65±0.23</b>	62.03±0.61	40.22±0.64	30.80±0.53	38.12±0.59	46.29±0.71
2	59.77±0.35	61.05±0.21	49.45 ± 0.16	58.22±0.14	67.82±0.31	<b>51.25±0.39</b>	32.90±0.24	41.31±0.26	50.64±0.38
3	61.40±0.22	61.55±0.14	51.33 ± 0.17	58.04±0.14	<b>70.60±0.31</b>	46.89±0.26	39.38±0.29	46.27±0.22	51.14±0.25
4	65.75±0.18	<b>63.40±0.19</b>	54.58 ± 0.14	<u>58.28±0.16</u>	<u>68.26±0.36</u>	41.49±0.23	42.42±0.16	49.61±0.20	58.02±0.25
5	68.15±0.25	62.35±0.19	55.66 ± 0.10	56.96±0.17	66.72±0.24	33.13±0.35	45.49±0.26	53.38±0.23	60.67±0.29
All	<b>69.12±0.21</b>	62.16±0.22	<b>56.79 ± 0.12</b>	57.21±0.17	64.02±0.21	41.64±0.44	<b>45.94±0.30</b>	<b>56.36±0.23</b>	<b>62.34±0.25</b>

Table 5: Single Chain results, Accuracy, Higher is better ↑

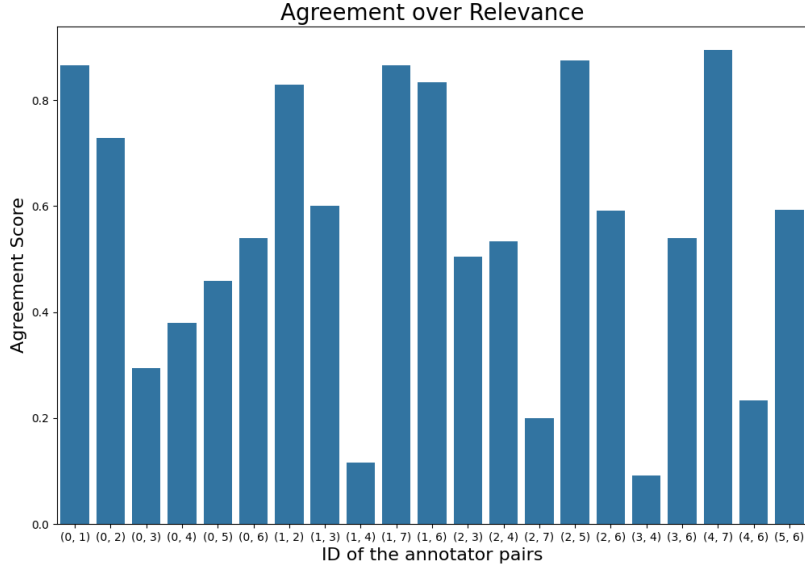


Figure 9: Annotator Agreement over the Relevance of a given reasoning step

Attribute	Overall Agreement
Relevance	0.55
Math Accuracy	0.81
Logical Consistency	0.53
Overall Correctness	0.66

Table 6: Overall Agreement by Attribute

Attribute	Overall Agreement
Relevance	0.71
Math Accuracy	0.84
Logical Consistency	0.75
Overall Correctness	0.81

Table 7: Overall Agreement by Attribute computed using Gwet’s AC1

and recorded the resulting final scores.<sup>9</sup>

We summarize our results in Table 10. We remark that even for modest correlations, the performance increase is over 20% relative, in line with what we observed empirically with real data.

## J Examples

We provide two qualitative examples in Figures 15, ??, and 17, comparing the solutions chosen by

<sup>9</sup>The correctness of reasoning chains in the artificially generated data is determined based on the sampled scores representing human judgments. A reasoning chain is considered correct if it consists of more than 75% correct reasoning steps.

the lowest perplexity method and by our proposed method, that of using the verifier scores.

1208

1209

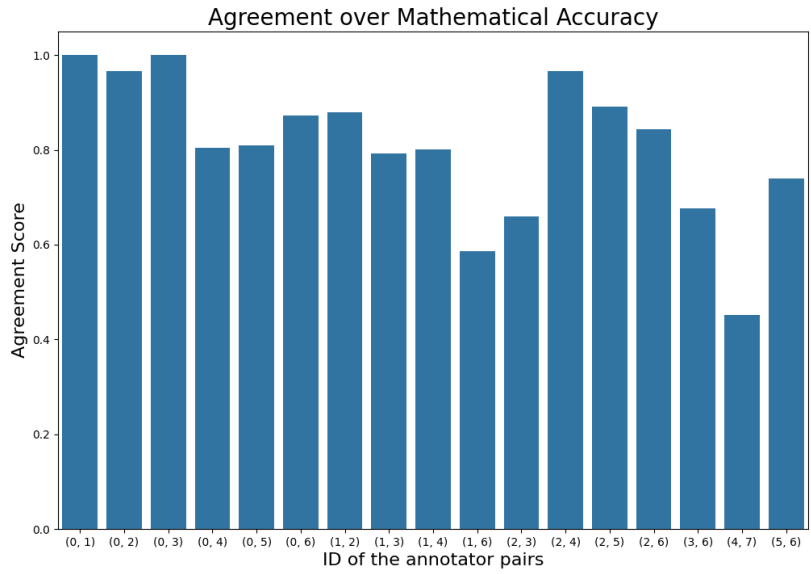


Figure 10: Annotator Agreement over the Mathematical Accuracy of a given reasoning step

Attribute	Overall Agreement
Relevance	0.83
Math Accuracy	0.89
Logical Consistency	0.86
Overall Correctness	0.89

Table 8: Overall Agreement by Attribute computed using the Naive approach

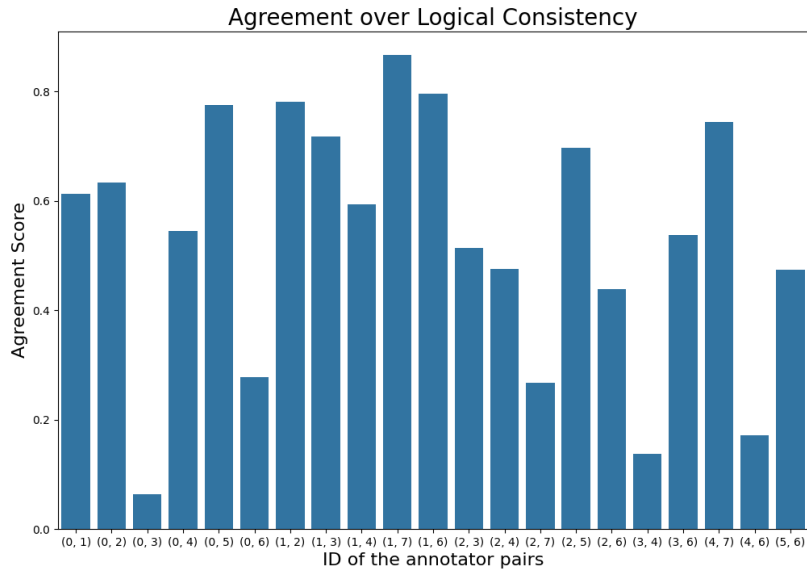


Figure 11: Annotator Agreement over the Logical Consistency of a given reasoning step

Metric\Annotator	A0	A1	A2	A3	A4	A5	A6	A7
Ratio of steps annotated as Irrelevant	0.41	0.37	0.24	0.33	0.17	0.22	0.36	0.35
Ratio of steps annotated as Mathematically incorrect	0.06	0.05	0.03	0.03	0.07	0.06	0.11	0.03
Ratio of steps annotated as Logically Inconsistent	0.4	0.35	0.25	0.31	0.28	0.23	0.28	0.34
Ratio of steps annotated as Overall Incorrect	0.43	0.38	0.31	0.37	0.31	0.24	0.34	0.36

Table 9: The ratio of steps each annotator annotated as Irrelevant, Mathematically incorrect, Logically Inconsistent, and Overall Incorrect.

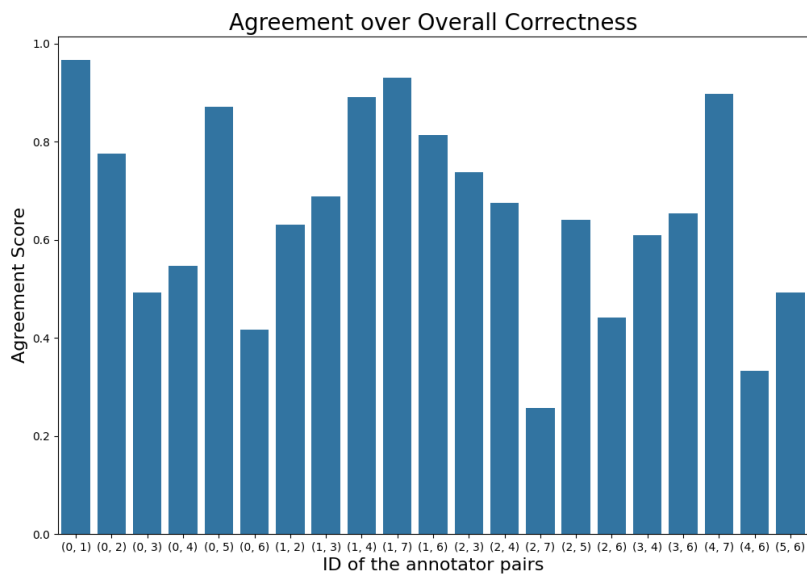
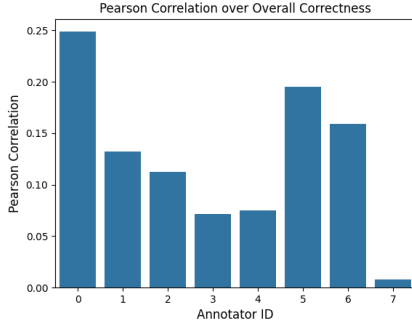
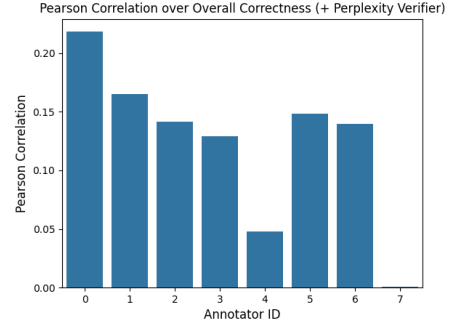


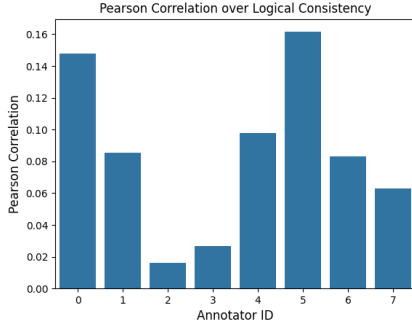
Figure 12: Annotator Agreement over the Overall Correctness of a given reasoning step



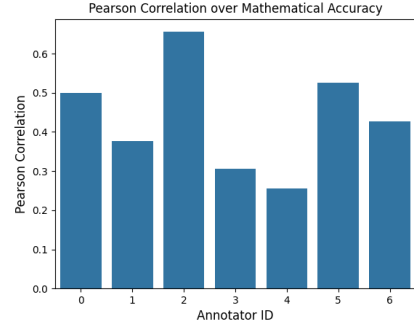
(a) Correlation score between our overall score (excluding the perplexity verifier) and the annotators assessment.



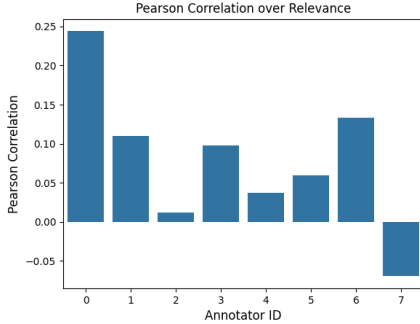
(b) Correlation score between our overall score (including the perplexity verifier) and the annotators assessment.



(c) Correlation score between our Logical Consistency Verifier and the annotators assessment.



(d) Correlation score between our Mathematical Accuracy Verifier and the annotators assessment.



(e) Correlation score between our Relevance Verifier and the annotators assessment.

Figure 13: Correlations between the scores of the verifiers and the human annotators' assessment. Additionally, we consider two scores for the *Overall Assessment* attribute. One was computed only with the three verifiers and the second one was computed with the three verifiers, together with the perplexity verifier.

Method	Score
Model with no verifier	0.19
Model with verifier with <i>correlation</i> = 0.075	0.22
Model with verifier with <i>correlation</i> = 0.1	0.23
Model with verifier with <i>correlation</i> = 0.15	0.26
Model with verifier with <i>correlation</i> = 0.25	0.31
Model with verifier with <i>correlation</i> = 0.5	0.48
Model with verifier with <i>correlation</i> = 0.75	0.72
Model with verifier with <i>correlation</i> = 1.0	0.98

Table 10: Scores of the proposed method given various level of correlations. Even for modest correlations, the performance increase is over 20% relative.

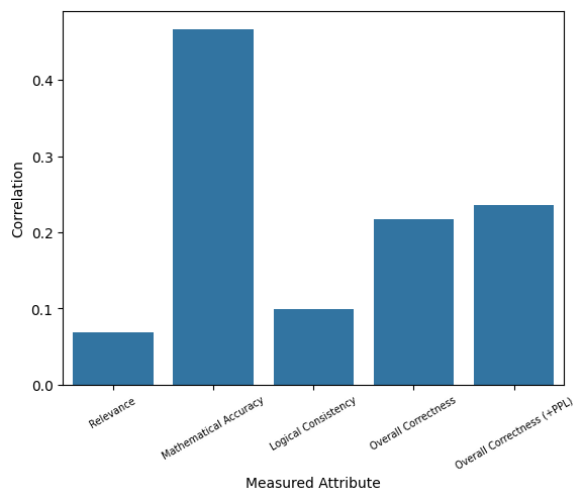


Figure 14: Correlation between the scores of our proposed verifiers and the human assessments. Additionally, we include the correlations between the human assessment of the overall correctness of a given reasoning step and the aggregated score (with and without perplexity).

**Problem Statement**

Antoine's french onion soup recipe calls for 2 pounds of onions. He likes to double that amount. His soup serves 6 people. The onions are currently on sale for \$2.00 a pound. He also needs 2 boxes of beef stock, that are also on sale for \$2.00 a box. What is the cost per serving? (Round to the nearest integer.)

**Gold Answer**

\$2

**Lowest perplexity solution**

- The recipe calls for 2 pounds of onions, so that's \$4.00 for the onions.
- Antoine doubles the amount of onions, so he uses 4 pounds for his recipe.
- That means he needs to buy 2 boxes of beef stock, which are on sale for \$2.00 each.
- The total cost for onions and beef stock is \$8.00.
- Therefore, the cost per serving is \$1.33 (rounded to the nearest integer).

**Solution according to verifiers**

- Antoine likes to double the amount of onions, so he uses 4 pounds of onions.
- Each onion costs \$2.00, so 4 pounds of onions cost  $4 * \$2.00 = \$\ll 4*2=8 \gg 8.00$ .
- He also needs 2 boxes of beef stock, which cost  $2 * \$2.00 = \$\ll 2*2=4 \gg 4.00$ .
- Therefore, the total cost for all ingredients is  $\$8.00 + \$4.00 = \$\ll 8+4=12 \gg 12.00$ .
- Since the soup serves 6 people, each person will eat  $\$12.00 / 6 = \$\ll 12/6=2 \gg 2$  worth of soup.
- Therefore, the cost per serving is \$2.00.

Figure 15: Qualitative example of the types of solutions preferred by our proposed method.

**Problem Statement**

Pam and Fred went to a carnival. Pam rode the roller coaster 2 times while Fred rode it 4 times. After that, each of them decided to ride the luge 2 times. If each ride cost 6 tickets, how many tickets did they use that day?

**Gold Answer**

60

**Lowest perplexity solution**

- Pam rode the roller coaster 2 times, so that's  $2 * 6 = 12$  tickets.
- Fred rode the roller coaster 4 times, so that's  $4 * 6 = 24$  tickets.
- After that, they each rode the luge 2 times, so that's  $2 * 2 * 6 = 24$  tickets.
- Finally, Pam rode the roller coaster 2 times, so that's  $2 * 6 = 12$  tickets.
- Fred rode the roller coaster 4 times, so that's  $4 * 6 = 24$  tickets.
- After that, they each rode the luge 2 times, so that's  $2 * 2 * 6 = 24$  tickets.
- Finally, Pam rode the roller coaster 2 times,

**Solution according to verifiers**

- Pam rode the roller coaster 2 times, which means she used  $2 * 6 = 12$  tickets.
- Fred rode the roller coaster 4 times, which means he used  $4 * 6 = 24$  tickets.
- After riding the roller coaster, Pam and Fred rode the luge 2 times each, which means they used  $2 * 2 * 6 = 24$  tickets in total.
- So, they used a total of  $12 + 24 + 24 = 60$  tickets that day.

Figure 16: Qualitative example of the types of solutions preferred by our proposed method.



**Problem Statement**

Eliza's rate per hour for the first 40 hours she works each week is \$10. She also receives an overtime pay of 1.2 times her regular hourly rate. If Eliza worked for 45 hours this week, how much are her earnings for this week?

**Gold Answer**

\$460

**Lowest perplexity solution**

- Eliza works for 40 hours at a rate of \$10 per hour, so her earnings for the first 40 hours are \$400.
- Since she worked 5 more hours than the 40 hours at her regular rate, she will receive 1.2 times her regular hourly rate for those 5 hours. (1.2 times 10 is 12)
- Therefore, Eliza's earnings for the 45 hours she worked this week are \$420. ( $\$400 + \$20$ )
- So, Eliza's earnings for this week are \$420.

**Solution according to verifiers**

- Eliza's regular hourly rate for the first 40 hours is \$10.
- Her overtime rate is 1.2 times the regular rate, which is \$12.
- She worked for 45 hours this week.
- So, for the first 40 hours, her earnings are \$400 ( $\$10/\text{hour} \times 40 \text{ hours}$ ), and for the 5 overtime hours, her earnings are \$60 ( $\$12/\text{hour} \times 5 \text{ hours}$ ).
- Adding these two amounts, Eliza's total earnings for this week are \$460 ( $\$400 + \$60$ ).

Figure 17: Qualitative example of the types of solutions preferred by our proposed method.